

Estimating Propensity Parameters using Google PageRank and Genetic Algorithms

David Murrugarra^{1*} Jacob Miller¹ Alex Mueller¹

¹Department of Mathematics, University of Kentucky, Lexington, KY 40506-0027 USA.

Abstract

Stochastic Boolean networks, or more generally, stochastic discrete networks, are an important class of computational models for molecular interaction networks. The stochasticity stems from the updating schedule. Standard updating schedules include the synchronous update, where all the nodes are updated at the same time, and the asynchronous update where a random node is updated at each time step. The former produces a deterministic dynamics while the latter a stochastic dynamics. A more general stochastic setting considers propensity parameters for updating each node. Stochastic Discrete Dynamical Systems (SDDS) are a modeling framework that considers two propensity parameters for updating each node and uses one when the update has a positive impact on the variable, that is, when the update causes the variable to increase its value, and uses the other when the update has a negative impact, that is, when the update causes it to decrease its value. This framework offers additional features for simulations but also adds a complexity in parameter estimation of the propensities. This paper presents a method for estimating the propensity parameters for SDDS. The method is based on adding noise to the system using the Google PageRank approach to make the system ergodic and thus guaranteeing the existence of a stationary distribution. Then with the use of a genetic algorithm, the propensity parameters are estimated. Approximation techniques that make the search algorithms efficient are also presented and Matlab/Octave code to test the algorithms are available at <http://www.ms.uky.edu/~dmu228/GeneticAlg/Code.html>.

Keywords:— Boolean Networks, Stochastic Systems, Propensity Parameters, Markov Chains, Google PageRank, Genetic Algorithms, Stationary Distribution

*Correspondence: murrugarra@uky.edu

1 Introduction

Mathematical modeling has been widely applied to the study of biological systems with the goal of understanding the important properties of the system and to derive useful predictions about the system. The type of systems of interest ranges from the molecular to ecological systems. At the cellular level, gene regulatory networks (GRN) have been extensively studied to understand the key mechanisms that are relevant for cell function. GRNs represent the intricate relationships among genes, proteins, and other substances that are responsible for the expression levels of mRNA and proteins. The amount of these gene products and their temporal patterns characterize specific cell states or phenotypes [24].

Gene expression is inherently stochastic with randomness in transcription and translation. This stochasticity is usually referred to as noise and it is one of the main drivers of variability [27]. Variability has an important role in cellular functions, and it can be beneficial as well as harmful [7, 14]. Modeling stochasticity is an important problem in systems biology. Different modeling approaches can be found in the literature. Mathematical models can be broadly divided into two classes: continuous, such as systems of differential equations and discrete, such as Boolean networks and their generalizations. This paper will focus on discrete stochastic methods. The Gillespie algorithm [8, 9] considers discrete states but continuous time. In this work, we will focus on models where the space as well as the time are discrete variables. For instance, Boolean networks (BNs) are a class of computational models in which genes can only be in one of two states: ON or OFF. BNs and, in general, multistate models, which allow genes to take on more than two states, have been effectively used to model biological systems such as the *p53-mdm2* system [1, 5, 25], the *lac* operon [37], the yeast cell cycle network [18], the Th regulatory network [21], *A. thaliana* [3], and many other systems [2, 6, 10, 11, 30, 39].

Stochasticity in Boolean networks has been studied in different ways. The earliest approach to introduce stochasticity into BNs was the asynchronous update, where a random node is updated at each time step [34]. Another approach considers update sequences that can change from step to step [22, 29]. More sophisticated approaches include Probabilistic Boolean Networks (PBNs) [31] and their variants [17, 19]. PBNs consider stochasticity at the function level where each node can use multiple functions with a switching probability from step to step. SDDS [25] is a simulation framework similar to PBNs but the key difference is how the transition probabilities are calculated. SDDS considers two propensity parameters for updating each node. These

parameters resemble the propensity probabilities in the Gillespie algorithm [8, 9]. This extension gives a more flexible simulation setup as a generative model but adds the complexity of parameter estimation of the propensity parameters. This paper provides a method for computing the propensity parameters for SDDS.

For completeness, in the following subsection, we will define the stochastic framework to be used in remainder of the paper.

Stochastic Framework

In this paper we will focus on the stochastic framework introduced in [25] referred to as Stochastic Discrete Dynamical Systems (SDDS). This framework is a natural extension of Boolean networks and is an appropriate setup to model the effect of intrinsic noise on network dynamics. Consider the discrete variables x_1, \dots, x_n that can take values in finite sets S_1, \dots, S_n , respectively. Let $S = S_1 \times \dots \times S_n$ be the Cartesian product. A *SDDS* in the variables x_1, \dots, x_n is a collection of n triplets

$$F = \{f_i, p_i^\uparrow, p_i^\downarrow\}_{i=1}^n$$

where

- $f_i : S \rightarrow S_i$ is the update function for x_i , for all $i = 1, \dots, n$.
- p_i^\uparrow is the activation propensity.
- p_i^\downarrow is the degradation propensity.
- $p_i^\uparrow, p_i^\downarrow \in [0, 1]$. These are the parameters of interest in this paper.

The stochasticity originates from the propensity parameters p_k^\uparrow and p_k^\downarrow , which should be interpreted as follows: If there would be an activation of x_k at the next time step, i.e., if $s_1, s_2 \in S_k$ with $s_1 < s_2$ and $x_k(t) = s_1$, and $f_k(x_1(t), \dots, x_n(t)) = s_2$, then $x_k(t+1) = s_2$ with probability p_k^\uparrow . The degradation probability p_k^\downarrow is defined similarly. SDDS can be represented as a Markov chain by specifying its transition matrix in the following way. For each variable x_i , $i = 1, \dots, n$, the probability of changing its value is given by

$$Prob(x_i \rightarrow f_i(x)) = \begin{cases} p_i^\uparrow, & \text{if } x_i < f_i(x), \\ p_i^\downarrow, & \text{if } x_i > f_i(x), \\ 1, & \text{if } x_i = f_i(x), \end{cases}$$

and the probability of maintaining its current value is given by

$$Prob(x_i \rightarrow x_i) = \begin{cases} 1 - p_i^\uparrow, & \text{if } x_i < f_i(x), \\ 1 - p_i^\downarrow, & \text{if } x_i > f_i(x), \\ 1, & \text{if } x_i = f_i(x). \end{cases}$$

Let $x, y \in S$. The transition from x to y is given by

$$a_{xy} = \prod_{i=1}^n Prob(x_i \rightarrow y_i). \quad (1)$$

Notice that $Prob(x_i \rightarrow y_i) = 0$ for all $y_i \notin \{x_i, f_i(x)\}$.

Then the transition matrix is given by

$$A = (a_{xy})_{x,y \in S} \quad (2)$$

The dynamics of SDDS depends on the transition probabilities a_{xy} , which depend on the propensity values and the update functions. Online software to test examples is available at <http://adam.plantsimlab.org/> (choose Discrete Dynamical Systems (SDDS) in the model type).

In Markov chain notation, the transition probability $a_{xy} = p(X_t = x | X_{t-1} = y)$ represents the probability of being in state x at time t given that system was in state y at time $t - 1$. If $\pi_t = p(X_t = x)$ represents the probability of being in state x at time t , then we will assume that π is a row vector containing the probabilities of being in state x at time t for all $x \in S$. If π_0 is the initial distribution at time $t = 0$, then at time $t = 1$,

$$\pi_1 = \sum_{x \in S} \pi_0(x) a_{xy}. \quad (3)$$

If we iterate Equation 3 and if we get to the point where

$$\pi = \sum_{x \in S} \pi(x) a_{xy} \quad (4)$$

then we will say that the Markov chain has reached a stationary distribution and that π is the stationary distribution.

2 Methods

In this section we describe a method for estimating the propensity parameters for SDDS. The approach is based on adding noise to the system using the

Google PageRank [4, 16, 23] strategy to make the system ergodic and thus guaranteeing the existence of a stationary distribution and then with the use of a genetic algorithm the propensity parameters are estimated. To guarantee the existence of a stationary distribution we use a special case of the Perron-Frobenius Theorem.

Theorem 2.1 (Perron-Frobenius). *If \mathbf{A} is a regular $m \times m$ transition matrix with $m \geq 2$, then*

- *For any initial probability vector π_0 , $\lim_{n \rightarrow \infty} \mathbf{A}^n \pi_0 = \pi$.*
- *The vector π is the unique probability vector which is an eigenvector of \mathbf{A} associated with the eigenvalue 1.*

A proof of Theorem 2.1 can be found in Chapter 10 of [16].

Theorem 2.1 ensures a unique stationary distribution π provided that we have a regular transition matrix, that is, if some power \mathbf{A}^k contains only strictly positive entries. However, the transition matrix \mathbf{A} of SDDS given in Equation 2 might not be regular. In the following subsection, we use a similar approach to the Google's PageRank algorithm to add noise to the system to obtain a new transition matrix that is regular.

2.1 PageRank Algorithm

For simplicity, consider a SDDS, $F = \{f_i, p_i^\uparrow, p_i^\downarrow\}_{i=1}^n$ where $f_i : S \rightarrow S_i$, $S = \mathbb{K}^n$, and $|\mathbb{K}| = p$. Then its transition matrix \mathbf{A} given in Equation 2 is a $p^n \times p^n$ matrix. To introduce noise into the system we consider the Google Matrix

$$\mathbf{G} = g\mathbf{A} + (1 - g)\mathbf{K}, \quad (5)$$

where g is a constant number in the interval $[0, 1]$ and \mathbf{K} is a $p^n \times p^n$ matrix all of whose columns are the vector $(1/p^n, \dots, 1/p^n)$. The matrix \mathbf{G} in Equation 5 is a regular matrix and then we can use Theorem 2.1 to get a stationary distribution for \mathbf{G} ,

$$\pi = \pi_{\mathbf{G}} = (\pi_1, \dots, \pi_{p^n}) \quad (6)$$

This stationary distribution reflects the dynamics of the SDDS $F = \{f_i, p_i^\uparrow, p_i^\downarrow\}_{i=1}^n$. The importance of a state $x \in S$ can be measured by the size of the corresponding entry π_x in the stationary distribution of Equation 6. For instance, for ranking the importance of the states in a Markov chain one can use the size of the corresponding entries in the stationary distribution. We will refer to this entry π_x as the PageRank score of x .

2.2 Genetic Algorithm

The entries of the stationary distribution π in Equation 6 can also be interpreted as occupation times for each state. Thus it gives the probability of being at a certain state. Now suppose that we start with a desired stationary distribution $\pi^* = (\pi_1^*, \dots, \pi_{p^n}^*)$. We have developed a genetic algorithm that initializes a population of random propensity matrices and searches for a propensity matrix $prop^*$ such that its stationary distribution $\pi = (\pi_1, \dots, \pi_{p^n})$ gets closer to the desired stationary distribution π^* . That is, we search for propensity matrices such that the distance between π and π^* is minimized,

$$\min_{p_i^\uparrow, p_i^\downarrow} d(\pi, \pi^*) \quad \text{or} \quad \min_{p_i^\uparrow, p_i^\downarrow} |\pi(j) - \pi^*(j)| \quad (7)$$

The pseudocode of this genetic algorithm is given in Algorithm 1 and it has been implemented in Octave/Matlab and our code can be downloaded from <http://www.ms.uky.edu/~dmu228/GeneticAlg/Code.html>.

Algorithm 1 Genetic Algorithm with PageRank.

Require: Functions: $F = (f_1, \dots, f_n)$, number of generations: $NumGen$, population size: $PopSize$, states of interest: $States$, and desired probabilities: $\pi^* = \pi^*(States)$.

Ensure: Propensity parameters: **prop***

```
1: procedure GENETICGOOGLE( $F, NumGen, PopSize, \pi^*$ )
2:    $PopPropensities \leftarrow$  initialize a population of propensity matrices.
3:    $[fitnesses, \min(PopPropensities)] = FITNESSGOOGLE(F, PopPropensities, \pi^*)$ 
4:   for  $i=1, \dots, NumGen$  do
5:      $NewPropensities \leftarrow$  initialize new population of propensities.
6:     for  $j=1, \dots, PopSize$  do
7:       if  $rand < fitnesses(j)$  then
8:          $parent1(j) = PopPropensities(j)$ 
9:       else
10:         $parent2(j) = PopPropensities(j)$ 
11:         $children = Crossover(parent1, parent2, mut, \sigma)$ 
12:         $NewPropensities(j) = children$ 
13:       $[fitnesses, \min(NewPropensities)] = FITNESSGOOGLE(F, NewPropensities, \pi^*)$ 
14:       $PopPropensities = NewPropensities$ .
15:    prop* =  $\min(NewPropensities)$ .
16: function FITNESSGOOGLE( $F, PopPropensities, \pi^*$ )
17:   for  $i = 1, \dots, \text{length}(PopPropensities)$  do    ▷ For each propensity
matrix.
18:      $\pi = \begin{cases} PageRank(F, PopPropensities(i)) & \text{for exact distribution, see Equation 6,} \\ ESTIMATESTADIST(F, c, NumIter, g) & \text{for estimated distribution, see Algorithm 2.} \end{cases}$ 
19:      $d = d(\pi, \pi^*)$  ▷ We used a weighted distance to give predominance
to important states.
20:      $fitnesses(i) = e^{(-d^2/s)}$ 
21:   return  $([fitnesses, \min(PopPropensities)])$  ▷ Keep propensity with
minimum fitness.
22: function Crossover( $parent1, parent2, mut, \sigma$ )
23:    $NewProp \leftarrow$  initialize new propensity matrix.
24:    $DivLine =$  random integer between 1 and  $\text{length}(parent1)$ .
25:   for  $i = 1, \dots, \text{length}(parent1)$  do
26:     if  $i < DivLine$  then
27:        $NewProp(i) = parent1(i)$ 
28:     else
29:        $NewProp(i) = parent2(i)$ 
30:     if  $rand < mut$  then    7
31:        $NewProp(i) = NewProp(i) + normrand(0, \sigma)$     ▷ Introduce
mutation.
32:   return ( $NewProp$ )
```

2.3 Estimating the stationary distribution

The genetic algorithm, Algorithm 1, uses the exact stationary distribution through PageRank (see Equation 6) which is computationally expensive for larger models. Here we present an efficient algorithm for estimating the stationary distribution based on a random walk. The expensive part of Algorithm 1 is the calculation of the stationary distribution π in Equation 6. We have implemented an algorithm for estimating the stationary distribution by doing a random walk using SDDS as a generative model; see Algorithm 2. The idea behind Algorithm 2 is to use SDDS for simulating from an initial state according to the transition probabilities given in Equation 5. That is, we initialize the simulation at an initial state $x \in S$ and then with probability g we move to another state $y \in S$ according to a_{xy} (see Equation 1) and with probability $1 - g$ we jump to a random node. We repeat this process for a given number of iterations. At the end of a maximum number of iterations, we count how often we visited each state and the normalized frequencies will be the approximated stationary distribution.

To make the genetic algorithm more efficient, we used the estimated stationary distribution described in the previous paragraph. Thus, within the fitness function of Algorithm 1, we use the estimated stationary distribution to assess the fitness of the generated propensity matrices. The pseudocode for this new algorithm is the same as Algorithm 1, the only change is in the fitness function. This version of the algorithm has also been implemented in Octave/Matlab and our code can be found in <http://www.ms.uky.edu/~dmu228/GeneticAlg/Code.html>.

Results

We test our methods using two published models that are appropriate for changing the stationary distribution under the choice of different propensity parameters. The first model is a Boolean network while the second is a multistate model. In both models bistability has been observed but the basin size of one of the attractors under the synchronous update is much larger than the basin of the other attractor, and thus the stationary distribution will be more concentrated in one of the attractors. We will use our methods to change the stationary distribution in favor of the attractor with a smaller basin.

Example 2.2. *Lac-operon network.* *The lac-operon in E. coli [12] is one of the best studied gene regulatory networks. This system is responsible for*

Algorithm 2 Estimate Stationary Distribution.

Require: Functions: $F = (f_1, \dots, f_n)$, propensities: c , number of iterations: $NumIter$, noise: g .

Ensure: Estimated stationary distribution π

```
 $\pi = \text{ESTIMATESTADIST}(F, c, NumIter, g)$ 
2: return  $\pi$ 
   function ESTIMATESTADIST( $F, c, NumIter, g$ )
4:    $distribution \leftarrow$  initialize frequency vector.
    $s \leftarrow$  initialize random initial state.
6:   for  $i=1, \dots, NumIter$  do
       if  $rand < g$  then
8:          $y =$  random state between 1 and  $p^n$ .
       else
10:         $y = SDDS.\text{nextstate}(s, c)$ 
         $distribution(y) = +1$  increase state frequency.
12:         $sum =$  total frequencies.
         $\pi = distribution/sum$ 
14:   return  $\pi$ 
```

the metabolism of lactose in the absence of glucose. This system exhibits bistability in the sense that the operon can be either ON or OFF, depending on the presence of the preferred energy source: glucose. A Boolean network for this system has been developed in [37]. This model considers the following 10 components

$$\begin{aligned} x_1 &= M: \text{lac mRNA}, & x_2 &= P: \text{lac permease}, \\ x_3 &= B: \text{lac}\beta\text{-galactosidase}, & x_4 &= C: \text{CAP}, \\ x_5 &= R: \text{repressor}, & x_6 &= Rm: \text{repressor at medium concentration}, \\ x_7 &= A: \text{allolactose}, & x_8 &= Am: \text{allolactose at medium concentration}, \\ x_9 &= L: \text{lactose}, & x_{10} &= Lm: \text{lactose at medium concentration}, \end{aligned} \tag{8}$$

and the Boolean rules are given by

$$\begin{aligned}
f_1 &= x_4 \wedge \overline{x_5} \wedge \overline{x_6}, \\
f_2 &= x_1, \\
f_3 &= x_1, \\
f_4 &= \overline{G_e}, \\
f_5 &= \overline{x_7} \wedge \overline{x_8}, \\
f_6 &= (\overline{x_7} \wedge \overline{x_8}) \vee x_5, \\
f_7 &= x_9 \wedge x_3, \\
f_8 &= x_9 \vee x_{10}, \\
f_9 &= x_2 \wedge L_e \wedge \overline{G_e}, \\
f_{10} &= ((L_{em} \wedge P) \vee L_e) \wedge \overline{G_e}.
\end{aligned} \tag{9}$$

where G_e , L_{em} , and L_e are parameters. G_e and L_e indicate the concentration of extracellular glucose and lactose, respectively. The parameter L_{em} indicates the medium concentration of extracellular lactose. For medium extracellular lactose, that is when $L_{em} = 1$ and $L_e = 0$, this system has two fixed points: $s_1 = (0, 0, 0, 1, 1, 1, 0, 0, 0, 0)$ and $s_2 = (1, 1, 1, 1, 0, 0, 0, 1, 0, 1)$ that represent the state of the operon being OFF and ON, respectively.

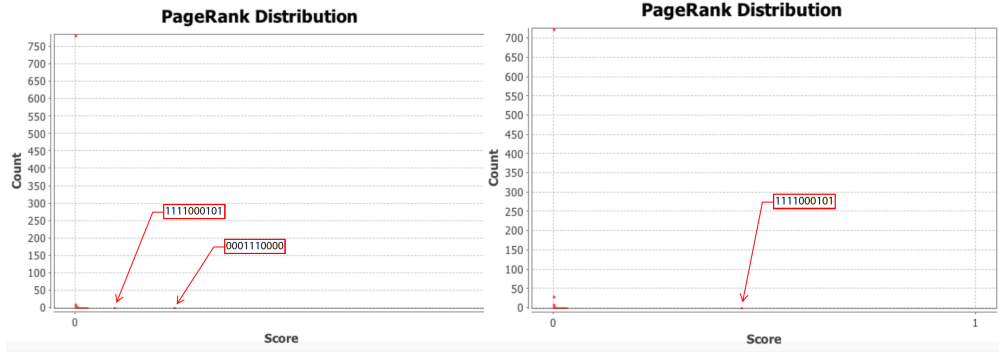
To test our method we calculated the stationary distribution of the system using Equation 6 with $g = 0.9$ in Equation 5. First we used the propensity values given in Equation 10 where all propensities are fixed to 0.9. This choice of parameters approximates the synchronous dynamics in the sense that each function has a 90% change of being used during the simulations and 10% chance of maintaining its current value. Under this selection of parameters, the fixed point s_1 has a PageRank score of 0.3346 while the other fixed point s_2 has a score of 0.0463, see Figure 1a and Table 1. Then we have applied our genetic algorithm to search for parameters that can increase the PageRank score of s_2 and decrease the score of s_1 . After doing so, we found the propensity parameters given in Equation 11. With this new set of parameters, the fixed point s_1 has a score of 0.0199 while s_2 has a score of 0.5485, see Figure 1b and Table 1. To appreciate the impact of the change in propensity parameters, we have plotted the state space of the system with both propensity matrices in Figure 2. The edges in blue in Figure 2 represent the most likely trajectory. Notice that in Figure 2b the trajectories are leading towards s_2 and the size of the labels of the nodes were scaled according to their PageRank score, see Table 1.

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
p_i^\uparrow	.9	.9	.9	.9	.9	.9	.9	.9	.9	.9
p_i^\downarrow	.9	.9	.9	.9	.9	.9	.9	.9	.9	.9

(10)

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
p_i^{\uparrow}	0.81	1.00	0.97	0.62	0.11	0.63	0.22	0.82	0.48	0.60
p_i^{\downarrow}	0.17	0.59	0.03	0.98	0.39	1.00	0.33	0.07	0.52	0.06

(11)



(a) Scores with propensities in Equation 10.

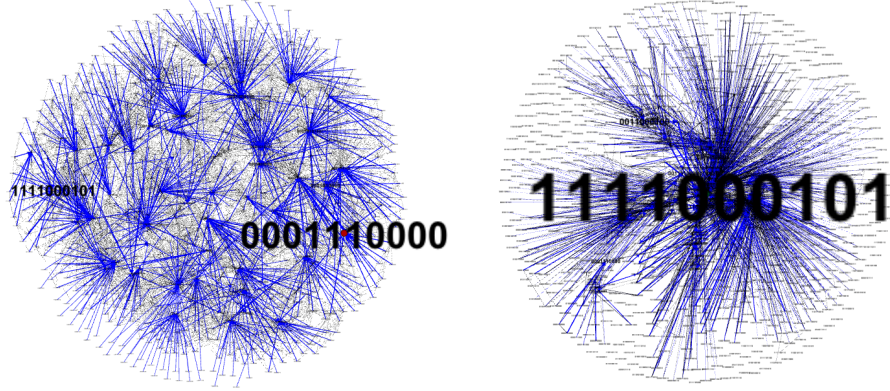
(b) Scores with propensities in Equation 11.

Figure 1: PageRank scores before and after the genetic algorithm. In each panel, the x -axis shows the PageRank scores while the y -axis shows the frequencies of states with the given scores in the x -axis (the exact scores for the states of interest are given in Table 1). Left panel shows the state space where all the propensities are equal to 0.9 while the right panel shows the state space where the propensity parameters were estimated using the genetic algorithm.

Propensities	Attractor	Score
In Equation 10 (all fixed to 0.9)	0001110000	0.3346
	1111000101	0.0463
In Equation 11 (genetic algorithm)	0001110000	0.0199
	1111000101	0.5485

Table 1: PageRank scores for the states of the attractors of the system. The order of variables in each vector state is $M, P, B, C, R, R_m, A, A_m, L, L_m$.

Example 2.3. Phage lambda infection. *The outcome of phage lambda infection is another system that has been widely studied over the last decades [13, 26, 32, 33, 38]. One of the earliest models that has been developed for this*



(a) State space with propensities in Equation 10. (b) State space with propensities in Equation 11.

Figure 2: State space comparison before and after the genetic algorithm. Left panel shows the state space where all the propensities are equal to 0.9 while the right panel shows the state space with the estimated propensity parameters using the genetic algorithm. The edges in blue represent the most likely trajectory. The size of the labels of the nodes were scaled according to their PageRank score.

system is the logical model by Thieffry and Thomas [33]. The regulatory genes considered in this model are: *CI*, *CRO*, *CII*, and *N*. Experimental reports [15, 28, 32, 33] have shown that, if the gene *CI* is fully expressed, all other genes are *OFF*. In the absence of *CRO* protein, *CI* is fully expressed (even in the absence of *N* and *CII*). *CI* is fully repressed provided that *CRO* is active and *CII* is absent.

This network is a bistable switch between lysis and lysogeny. Lysis is the state where the phage will be replicated, killing the host. Otherwise, the network will transition to a state called lysogeny where the phage will incorporate its DNA into the bacterium and become dormant. These cell fate differences have been attributed to the spontaneous changes in the timing of individual biochemical reaction events [20, 33].

In the model of Thieffry and Thomas [33], the first variable, *CI*, has three levels $\{0, 1, 2\}$, the second variable, *CRO*, has four levels $\{0, 1, 2, 3\}$, and the third and fourth variables, *CII* and *N*, are Boolean. Since the nodes of this model have different number of states, in order to apply our methods, we have extended the model so that all nodes have the same number of states. We have used the method given in [36] to extend the number of states such that all nodes have 5 states (the method for extending requires a prime number

for number of states so we have chosen 5 states). The method for extending the number of states preserves the original attractors. The update rules for this model are available with our code that is freely available. The extended model has a steady state, 2000, and a 2-cycle involving 0200 and 0300. The steady state 2000 represents lysogeny where *CI* is fully expressed while the other genes are *OFF*. The cycle between 0200 and 0300 represents lysis where *CRO* is active and other genes are repressed.

To test our method we calculated the stationary distribution of the system using Equation 6 with $g = 0.9$ in Equation 5. First we used the propensity values given in Equation 12 where all propensities are fixed to 0.9. This choice of parameters approximates the synchronous dynamics in the sense that each function has a 90% change of being used during the simulations and 10% chance of maintaining its current value. Under this selection of parameters, the fixed point 2000 has a PageRank score of 0.2772 while the states of the cycle 0200 and 0300 have scores of 0.2185 and 0.2108, respectively. Notice that this cycle attractor will have an overall score of 0.4293, see Figure 3a and Table 2. Then we have applied our genetic algorithm to search for parameters that can increase the PageRank score of the fixed point 2000 and found the propensity parameters given in Equation 13. With this new set of parameters, the fixed point 2000 has a score of 0.6040 while the states of the cycle 0200 and 0300 have scores of 0.0716 and 0.00016, respectively, see Figure 3b and Table 2. To appreciate the impact of the change in propensity parameters, we have plotted the state space of the system with both propensity matrices in Figure 4. The edges in blue in Figure 4 represent the most likely trajectory. Notice that in Figure 4b the trajectories are leading towards 2000 and the size of the labels of the nodes were scaled according to their PageRank score, see Table 2.

	<i>CI</i>	<i>CRO</i>	<i>CII</i>	<i>N</i>
p_i^\uparrow	.9	.9	.9	.9
p_i^\downarrow	.9	.9	.9	.9

(12)

	<i>CI</i>	<i>CRO</i>	<i>CII</i>	<i>N</i>
p_i^\uparrow	1.0000	0	0.4277	0.7968
p_i^\downarrow	0.3962	1.0000	0.6063	0.6946

(13)

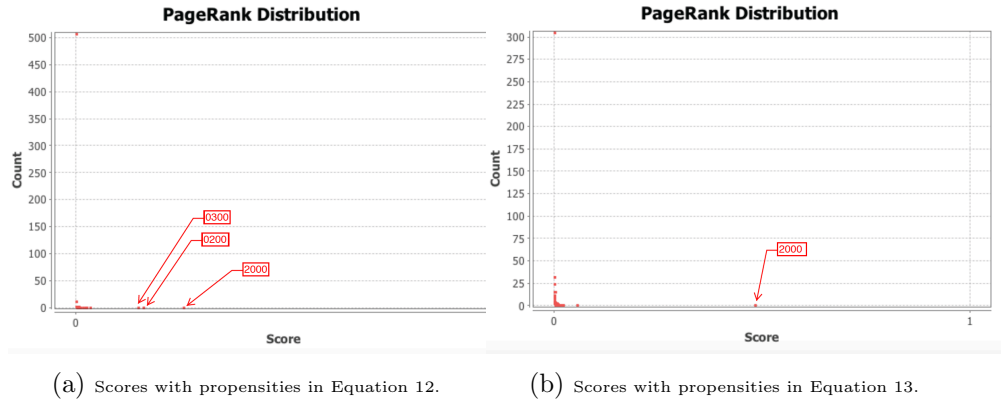


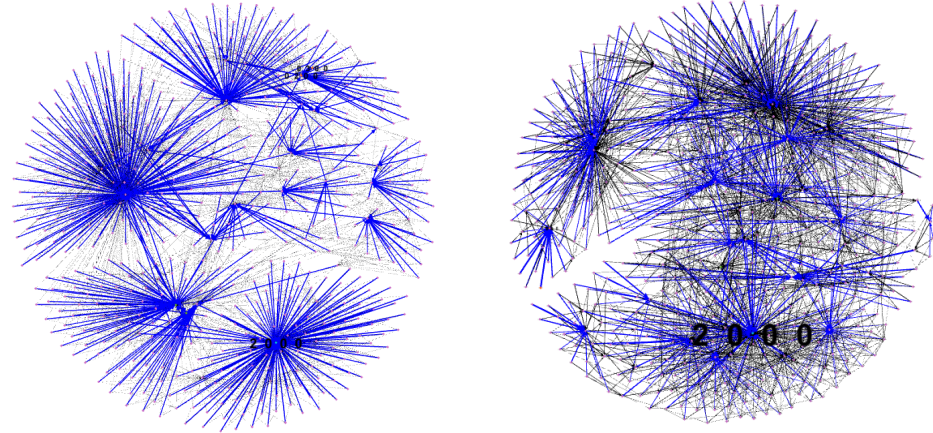
Figure 3: PageRank scores before and after the genetic algorithm. In each panel, the x -axis shows the PageRank scores while the y -axis shows the frequencies of states with the given scores in the x -axis (the exact scores for the states of interest are given in Table 2). Left panel shows the PageRank scores where all the propensities were equal to 0.9 while the right panel shows the scores where the propensity parameters were estimated using the genetic algorithm. The score for the state 2000 is 0.6040.

Propensities	Attractor	Score
In Equation 12 (all fixed to 0.9)	2000	0.2772
	0200	0.2185
	0300	0.2108
In Equation 13 (genetic algorithm)	2000	0.6040
	0200	0.0716
	0300	0.00016

Table 2: PageRank scores for the states of the attractors of the system. The order of variables in each vector state is CI, CRO, CII, N .

3 Discussion

Parameter estimation for stochastic models of biological networks is in general a very hard problem. This paper focuses on a class of stochastic discrete models, which is an extension of Boolean networks. The methods presented here use a well established approach for introducing noise into a system in a way that ergodicity and thus the existence of a unique stationary distribution is guaranteed. Then a genetic algorithm for calculating a set of parameters able to approximate a desired stationary distribution was developed. Also,



(a) State space with propensities in Equation 12. (b) State space with propensities in Equation 13.

Figure 4: State space comparison before and after the genetic algorithm. Left panel shows the state space where all the propensities are equal to 0.9 while the right panel shows the state space with the estimated propensity parameters using the genetic algorithm. The edges in blue represent the most likely trajectory. The size of the labels of the node were scaled according to their PageRank score.

techniques for approximating the stationary distribution at each iteration of the genetic algorithm that make the search process more efficient was applied.

One shortcoming of the method is that it is a stochastic method. That is, each time we run the algorithm we get a different result. An exhaustive investigation about the variance of the results is still missing. For the examples that we presented in the results section, the output of the algorithm might vary in about 20% of the reported propensities.

In parameter estimation, sometimes, it is useful to identify a smaller set of key parameters to estimate. This problem is out of the scope of the paper. However, for Boolean network models, one way to address this problem could be by using the different network reduction algorithms, for instance see [29, 35], to identify a smaller “core” network that preserves the important features of the dynamics of the original network. And then one could apply the parameter estimation techniques that are described in this paper. This type of approach could be especially useful if dealing with very large networks where running the genetic algorithm is computationally expensive.

4 Conclusions

In this paper we present an efficient method for estimating the parameters of a stochastic framework. The modeling framework is an extension of Boolean networks that uses propensity parameters for activation and inhibition. Parameter estimation techniques are needed whenever one needs to tune the propensity parameters of the stochastic system to reproduce a desired stationary distribution. For instance, if dealing with a bistable system and if it is desired to have the stationary distribution that have the PageRank score concentrated in one of the attractors of the system, then one needs to estimate the propensity parameters that represent such a desired distribution. Parameter estimation methods for this purpose were not available. In this paper we present a method for estimating propensity parameters given a desired stationary distribution for the system. We tested the method in one Boolean network with 10 nodes (where the size of the state space is $2^{10} = 1024$) and a multistate network with 4 nodes where each node has 5 states (where the size of the state space is $4^5 = 625$). For each system, we were able to redirect the system towards the attractor with the smaller PageRank score. The method is efficient and for the examples we have shown it can be run in few seconds in a laptop computer. Our code is available at <http://www.ms.uky.edu/~dmu228/GeneticAlg/Code.html>.

Conflict of Interest Statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Author Contributions

DM, JM, and AM designed the project. JM developed the genetic algorithm and implemented the Matlab code. AM run simulations for the results section. DM wrote the paper. All authors approved the final version of the manuscript.

Funding

DM was partially funded by a startup fund from the College of Arts and Sciences at the University of Kentucky.

Acknowledgments

The authors thank the reviewers for their insightful comments that have improved the manuscript.

References

- [1] Wassim Abou-Jaoudé, Djomangan A Ouattara, and Marcelle Kaufman. From structure to dynamics: frequency tuning in the p53-mdm2 network i. logical approach. *J Theor Biol*, 258(4):561–77, Jun 2009.
- [2] Réka Albert and Hans G Othmer. The topology of the regulatory interactions predicts the expression pattern of the segment polarity genes in drosophila melanogaster. *J Theor Biol*, 223(1):1–18, Jul 2003.
- [3] Enrique Balleza, Elena R Alvarez-Buylla, Alvaro Chaos, Stuart Kauffman, Ilya Shmulevich, and Maximino Aldana. Critical dynamics in genetic regulatory networks: examples from four kingdoms. *PLoS One*, 3(6):e2456, 2008.
- [4] Sergey Brin and Lawrence Page. Reprint of: The anatomy of a large-scale hypertextual web search engine. *Computer Networks*, 56(18):3825 – 3833, 2012. The {WEB} we live in.
- [5] Minsoo Choi, Jue Shi, Sung Hoon Jung, Xi Chen, and Kwang-Hyun Cho. Attractor landscape analysis reveals feedback loops in the p53 network that control the cellular response to dna damage. *Sci. Signal.*, 5(251):ra83, 2012.
- [6] Maria I Davidich and Stefan Bornholdt. Boolean network model predicts cell cycle sequence of fission yeast. *PLoS One*, 3(2):e1672, 2008.
- [7] Avigdor Eldar and Michael B Elowitz. Functional roles for noise in genetic circuits. *Nature*, 467(7312):167–73, Sep 2010.
- [8] Daniel T. Gillespie. Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical Chemistry*, 81(25):2340–2361, 1977.
- [9] Daniel T. Gillespie. Stochastic simulation of chemical kinetics. *Annual Review of Physical Chemistry*, 58(1):35–55, 2007. PMID: 17037977.
- [10] Tomáš Helikar, Naomi Kochi, Bryan Kowal, Manjari Dimri, Mayumi Naramura, Srikumar M Raja, Vimla Band, Hamid Band, and Jim A

- Rogers. A comprehensive, multi-scale dynamical model of erbb receptor signal transduction in human mammary epithelial cells. *PLoS One*, 8(4):e61757, 2013.
- [11] Tomás Helikar, John Konvalina, Jack Heidel, and Jim A Rogers. Emergent decision-making in biological signal transduction networks. *Proc Natl Acad Sci U S A*, 105(6):1913–8, Feb 2008.
 - [12] François Jacob and Jacques Monod. Genetic regulatory mechanisms in the synthesis of proteins. *Journal of Molecular Biology*, 3(3):318 – 356, 1961.
 - [13] Richard I. Joh and Joshua S. Weitz. To lyse or not to lyse: Transient-mediated stochastic fate determination in cells infected by bacteriophages. *PLoS Computational Biology*, 7(3), 2011.
 - [14] Mads Kaern, Timothy C Elston, William J Blake, and James J Collins. Stochasticity in gene expression: from theories to phenotypes. *Nat Rev Genet*, 6(6):451–64, Jun 2005.
 - [15] P Kourilsky. Lysogenization by bacteriophage lambda. i. multiple infection and the lysogenic response. *Mol Gen Genet*, 122(2):183–95, Apr 1973.
 - [16] David C. Lay. *Linear Algebra And Its Applications*. PEARSON, fourth edition, 2012.
 - [17] Ritwik Layek, Aniruddha Datta, Ranadip Pal, and Edward R Dougherty. Adaptive intervention in probabilistic boolean networks. *Bioinformatics*, 25(16):2042–8, Aug 2009.
 - [18] Fangting Li, Tao Long, Ying Lu, Qi Ouyang, and Chao Tang. The yeast cell-cycle network is robustly designed. *Proc Natl Acad Sci U S A*, 101(14):4781–6, Apr 2004.
 - [19] Jinghang Liang and Jie Han. Stochastic boolean networks: an efficient approach to modeling gene regulatory networks. *BMC Syst Biol*, 6:113, 2012.
 - [20] Harley H McAdams and Adam Arkin. Stochastic mechanism in gene expression. *Proceedings of the National Academy of Sciences*, 94(3):814–819, 1997.

- [21] Luis Mendoza. A network model for the control of the differentiation process in th cells. *Biosystems*, 84(2):101–14, May 2006.
- [22] Henning S. Mortveit and Christian M. Reidys. *An Introduction to Sequential Dynamical Systems*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2007.
- [23] Kevin P Murphy. *Machine learning: a probabilistic perspective*. The MIT Press, 2012.
- [24] David Murrugarra and Elena S Dimitrova. Molecular network control through boolean canalization. *EURASIP J Bioinform Syst Biol*, 2015(1):9, Dec 2015.
- [25] David Murrugarra, Alan Veliz-Cuba, Boris Aguilar, Seda Arat, and Reinhard Laubenbacher. Modeling stochasticity and variability in gene regulatory networks. *EURASIP Journal on Bioinformatics and Systems Biology*, 2012(1):5, 2012.
- [26] Mark Ptashne and A Genetic Switch. Phage lambda and higher organisms. *Cell & Blackwell Scientific, Cambridge, MA*, 1992.
- [27] Arjun Raj and Alexander van Oudenaarden. Nature, nurture, or chance: Stochastic gene expression and its consequences. *Cell*, 135(2):216–226, 2016/07/07 2008.
- [28] Louis Reichardt and AD Kaiser. Control of λ repressor synthesis. *Proceedings of the National Academy of Sciences*, 68(9):2185–2189, 1971.
- [29] Assieh Saadatpour, István Albert, and Réka Albert. Attractor analysis of asynchronous boolean models of signal transduction networks. *J Theor Biol*, 266(4):641–56, Oct 2010.
- [30] Assieh Saadatpour, Rui-Sheng Wang, Aijun Liao, Xin Liu, Thomas P Loughran, István Albert, and Réka Albert. Dynamical and structural analysis of a t cell survival network identifies novel candidate therapeutic targets for large granular lymphocyte leukemia. *PLoS Comput Biol*, 7(11):e1002267, Nov 2011.
- [31] Ilya Shmulevich, Edward R. Dougherty, Seungchan Kim, and Wei Zhang. Probabilistic boolean networks: a rule-based uncertainty model for gene regulatory networks. *Bioinformatics*, 18(2):261–274, 2002.

- [32] François St-Pierre and Drew Endy. Determination of cell fate selection during phage lambda infection. *Proc Natl Acad Sci U S A*, 105(52):20705–10, Dec 2008.
- [33] D Thieffry and R Thomas. Dynamical behaviour of biological regulatory networks—ii. immunity control in bacteriophage lambda. *Bull Math Biol*, 57(2):277–97, Mar 1995.
- [34] René Thomas and Richard D’Ari. *Biological feedback*. CRC Press, Boca Raton, 1990.
- [35] Alan Veliz-Cuba, Boris Aguilar, Franziska Hinkelmann, and Reinhard Laubenbacher. Steady state analysis of boolean molecular network models via model reduction and computational algebra. *BMC Bioinformatics*, 15:221, 2014.
- [36] Alan Veliz-Cuba, Abdul Salam Jarrah, and Reinhard Laubenbacher. Polynomial algebra of discrete models in systems biology. *Bioinformatics*, 26(13):1637–43, Jul 2010.
- [37] Alan Veliz-Cuba and Brandilyn Stigler. Boolean models can explain bistability in the *lac* operon. *Journal of Computational Biology*, 18(6):783–794, 2011.
- [38] Lanying Zeng, Samuel O Skinner, Chenghang Zong, Jean Sippy, Michael Feiss, and Ido Golding. Decision making at a subcellular level determines the outcome of bacteriophage infection. *Cell*, 141(4):682–691, May 2010.
- [39] Ranran Zhang, Mithun Vinod Shah, Jun Yang, Susan B Nyland, Xin Liu, Jong K Yun, Réka Albert, and Thomas P Loughran, Jr. Network model of survival signaling in large granular lymphocyte leukemia. *Proc Natl Acad Sci U S A*, 105(42):16308–13, Oct 2008.